

# Főnévi események automatikus detektálása függőségi elemző és WordNet alkalmazásával magyar nyelvű szövegeken

Subecz Zoltán<sup>1</sup>

<sup>1</sup>Pallasz Athéné Egyetem, Kecskemét  
subecz@szolf.hu

**Kivonat.** A természetes szövegekből történő információkinyerés egyik fontos részterülete a névelemek azonosítása mellett az események detektálása. Szövegekben lévő események detektálása és analizálása fontos szerepet tölt be számos számítógépes nyelvészeti alkalmazásban, mint például a kivonatolás és a válaszkérés. A szövegekben a legtöbb esemény igékhez kapcsolódik, és az igék általában eseményeket jelölnek. De az igéken kívül lehetnek események más szófajú szavak is pl. főnevek, igenevek stb. Munkánkban a szövegekben megtalálható főnévi események detektálásával foglalkoztunk. Jelen tanulmányunkban bemutatjuk gazdag jellemzőtérre alapuló gépi tanuló megközelítésünket, amely automatikusan képes magyar nyelvű szövegekben főnévi események detektálására függőségi elemző és WordNet alkalmazásával. A jellemzők mellé kiegészítő módszereket is alkalmaztunk, amelyek javították az eredményeket és a futási időt. Algoritmusainkat tesztadatbázisokon kiértékelve versenyképes eredményeket érnek el az eddig bemutatott angol és más nyelvű eredményekkel összehasonlítva.

**Kulcsszavak:** információkinyerés, eseménydetektálás, főnévi események detektálása, WordNet, függőségi elemzés

## 1 Bevezetés

A természetes szövegekből történő információkinyerés egyik fontos részterülete a névelemek azonosítása mellett az események detektálása [7]. Szövegekben lévő események detektálása és analizálása fontos szerepet tölt be számos számítógépes nyelvészeti alkalmazásban, mint például a kivonatolás és a válaszkérés. A szövegekben lévő események felismerése, analizálása, és hogy hogyan viszonyulnak egymáshoz időben, fontos a szöveg tartalmának megismerésében.

Az esemény, ami történik egy adott helyen és időben. A szövegekben a legtöbb esemény igékhez kapcsolódik, és az igék általában eseményeket jelölnek. De az igéken kívül lehetnek események más szófajú szavak is pl. főnevek, igenevek stb. Munkánkban a szövegekben megtalálható főnévi események detektálásával foglalkoztunk. Vannak olyan szavak (pl. írás), amelyek egyes mondatokban események, másokban pedig nem, ezért a szavak szöveggörnyezetét is elemezni kell. Jelen tanulmányunkban

bemutatjuk gazdag jellemzőtérre alapuló gépi tanuló megközelítésünket, amely automatikusan képes magyar nyelvű szövegekben főnévi események detektálására függőségi elemző és WordNet alkalmazásával. A rendszer bemenete egy token-szinten címkézett tanító korpusz. Modellünk jelöltjei a mondatok főnevei voltak.

A feladatokhoz gazdag jellemzőkészletre alapuló osztályozót használtunk. A jellemzők mellé kiegészítő módszereket is alkalmaztunk, amelyek javították az eredményeket és a futási időt. Módszerünket a Szeged Korpusz öt különböző doménjén vizsgáltuk meg.

Angol nyelvű szövegekre általában *konstituensfa alapú* szintaktikai elemzőt használunk az elő-feldolgozásnál az angol nyelv erősen konfiguratív tulajdonsága miatt, ahol is a legtöbb mondatszintű szintaktikai információt a szórenddel fejeznek ki. Ezzel ellentétben a magyar nyelv gazdag morfológiával és szabad szórenddel rendelkezik. A *függőségi fákkal* dolgozó elemzők különösen jól használhatóak szabad szórendű nyelvek elemzésére, így a magyarra is. Ezek ugyanis könnyebben teszik lehetővé az egymással nem szomszédos, de összetartozó szavak összekapcsolását is. Ezért mi a magyar nyelvű szövegeinkre *függőségi fákkal dolgozó elemzőt* használtunk.

Megoldásunkban a vizsgált szavak szemantikai jellemzéséhez felhasználtuk a magyar *WordNet*-et [10]. Mivel egy szóalakhoz több jelentés is tartozhat a *WordNet*-ben, ezért az egyes jelentések között egyértelműsítést végeztünk a *Lesk algoritmus*sal [8].

Algoritmusainkat tesztadatbázisokon kiértékelve, versenyképes eredményeket érnek el az eddig bemutatott angol és más nyelvű eredményekkel összehasonlítva.

## 2 Kapcsolódó munkák

Az EVITA [13] volt az első esemény felismerő eszközök egyike. Az eseményeket nyelvészeti és statisztikai technikák kombinálásának segítségével ismeri fel. Nyelvészeti ismereteken alapuló szabályokat használ fő jellemzőként. A főnévi esemény felismeréshez *WordNet* osztályokat is használ, valamint a főnevek szemantikai egyértelműsítésére Bayes osztályozót alkalmaz.

Boguraev és társa [2] gépi tanuláson alapuló módszert mutattak be automatikus esemény-annotáláshoz. A feladatot osztályozásra visszavezetve, RRM (robust risk minimization) osztályozót alkalmaztak. Jellemzőkként lexikai, morfológiai és szintaktikai chunk típusokat használtak két- és háromelemű ablakokban vizsgálva.

Bethard és társa [1] esemény felismerésre fejlesztették a STEP rendszert. Szintaktikai és szemantikai jellemzőket alkalmaztak és az esemény felismerési feladatot osztályozásként oldották meg. Gazdag jellemző készletet építettek be: lexikai, morfológiai, szintaktikai függőségi és választott *WordNet* osztályokat. E jellemzőkre alapozva Support Vector Machine (SVM) modellt implementáltak.

Llorens és társa [9] bemutattak egy kiértékelést esemény felismerésre. Szemantikai szerepeket adtak meg jellemzőként és CRF (Conditional Random Field) modellt építettek események felismeréséhez.

Jeong és társa [6] függőségi elemzőt használtak, de csak a jelölt főnév és a közvetlenül ahhoz kapcsolódó ige közötti kapcsolatot vizsgálták. Kombinált jellemzőket építettek be, az ige és a hozzá tartozó kapcsolat-típus párokat alkalmazva. A

WordNetet is használták, de jelentéségyértelműsítés nélkül. A MaxEnt osztályozási algoritmust a következő jellemzőkészlettel implementálták: felszíni, lexikai, szemantikai, függőségi alapú jellemzők. A jellemzőket a Kullback-Leibler divergencia módszerrel súlyozták.

Olasz szövegekre Caselli és társa [3] csak igéből képzett főnévi eseményekkel foglalkoztak, amihez a Weka döntési fa alapú osztályozót használták. Vizsgálták a jelölt argumentum struktúráját, az aspektuális módosítókat, a jelölt előtti és utáni 3-3 szófaját.

Spain szövegekre Peris és társa [11] szintén csak igéből képzett eseményekkel foglalkoztak. Osztályozásra a Weka döntési fa osztályozóját alkalmazták és külső főnévi lexikont használtak fel. Függőségi elemzőt alkalmaztak, de csak a jelölt főnév és a közvetlenül ahhoz kapcsolódó ige közötti kapcsolatot vizsgálták. Felhasználták a jelölt argumentum struktúráját is.

Német nyelvű szövegekre Gorzitze és társa [5] bootstrapping módszert alkalmaztak események felismerésre. Idővel kapcsolatos kifejezéseket és aspektuális igeiket kerestek a jelölt közelében. A jelölt és a közvetlenül ahhoz kapcsolódó ige kapcsolatát vizsgálták és szabály alapú függőségi elemzőt használtak.

Magyar szövegekre Subecz [14] detektált eseményeket, de csak igei és főnévi ige-névi eseményekkel foglalkozott. A következő jellemző készletet használta: felszíni, lexikai, morfológiai, szintaktikai, WordNet és frekvencia jellemzők. Ezen jellemzők mellett szabály alapú módszereket is alkalmazott.

### 3 Események, főnévi események

A szövegekben a legtöbb esemény ige-közvetlenül kapcsolódik, és az ige általános eseményeket jelölnek. De az ige-közvetlenül kívül lehetnek események más szó-fajájú szavak is például főnevek, ige-névek. Munkánkban a szövegekben megtalálható főnévi események detektálásával foglalkoztunk. Példák főnévi eseményekre: futás, építés, írás, háború, ünnepség.

A főnévi eseményeknek két nagy csoportja van: igéből képzettek (deverbális) és nem igéből képzettek (nem deverbális). Példa igéből képzett eseményekre: futás, írás. Példa nem igéből képzett eseményre: háború. Az igéből képzett főnevek két fő-fajtája az események és az eredmények, amelyeknél gyakori a kétértelműség is. Vannak olyan szavak (például írás), amelyek egyes mondatokban események, másokban pedig eredmények.

Például az *írás* főnév a következő mondatban esemény: *Az írás 5 órakor kezdődött.* Viszont a következő mondatban nem esemény, hanem eredmény: *Eloolvastuk az írást.* A többértelműség miatt nem elég a szóalak vizsgálata, a szövegkörnyezetet is elemezni kell.

### 4 Környezet

Alkalmazásunkban a Szeged Korpusz [4] egy részét használtuk fel a következő területekről: *üzleti rövidhírek, szépirodalom-fogalmazás, számítógépes szövegek, újsághí-*

*rek, jogi szövegek.* Tanításhoz és kiértékeléshez tízszeres keresztvalidációt alkalmaztunk. A mondatokat két nyelvész annotálta, az annotátorok közötti egyetértés Kappa = 0,7 volt.

A feladatokat *bináris osztályozásra* vezettük vissza. Az osztályozáshoz a *Weka* programcsomagnak<sup>1</sup> a J48-as döntési fa elemzőjét használtuk fel. A feladathoz alkalmaztuk még a Magyarlanc 2.0 programcsomagot is [16]. A csomagot magyar szövegek mondatokra és szavakra bontására, a szavak morfológiai elemzésére, majd szófaji egyértelműsítésére, és mondatok függőségi nyelvtan szerinti szintaktikai elemzésére alkalmaztuk. A Magyarlanc programcsomag is készít a szavakhoz morfológiai elemzést, de a HunMorph[15] elemzőcsomag sok esetben részletesebb elemzést ad, ezért ezt is felhasználtuk. Így a feladathoz két morfológiai elemzőt is alkalmaztunk.

Ahogy láttuk a kapcsolódó munkáknál, mások is használták függőségi elemzést. De az elemzőfában mindenki csak a jelölt és a vele közvetlen kapcsolatban lévő szavakat vizsgálta. Mi vizsgáltuk a jelölt és a fában tőle távolabbi igék kapcsolatát is. Modellünk jelöltjei a mondatok főnevei voltak. Számos esetben a jelölt alá egy részfa tartozott a függőségi fában.

A vizsgált főnevek szemantikai jellemzéséhez a magyar *WordNet*-et[10] alkalmaztuk. A WordNet hiperním hierarchiájában található szemantikai kapcsolatokat használtuk fel.

A szavak közötti szintaktikai kapcsolatok alapján a mondatok egy *függőségi fát* alkotnak. A fa legfelső eleme a *Root*. A fa *csomópontjaiban* vannak a mondat szavai, az *ágak* a szavak közötti *szintaktikai kapcsolatokat* reprezentálják. Ha a jelölt több szóból áll, akkor ezek a szavak egy részfát alkotnak a mondat fáján belül. A részfa a kiemelt szaván (fejszó, headword) keresztül kapcsolódik a fa többi részéhez.

*Statisztikai adatok:* A vizsgált korpusz 10000 mondatot tartalmaz. Jelöltek száma (főnevek): 48388 db. Pozitív jelöltek száma (esemény főnevek): 7626 db. A jelölteket két fő részre osztottuk a hasonló tulajdonságok alapján. Az egyik csoportba az igéből képzett főnevek, a másikba a többi főnév került. Igéből képzett jelöltek: 5325 db. Igéből képzett pozitív jelöltek: 4169 db. Nem igéből képzett jelöltek: 43063 db. Nem igéből képzett pozitív jelöltek: 3457 db.

## 5 Az osztályozás bemutatása

Az osztályozáshoz bináris osztályozót használtunk. A mondatok főnevei voltak a jelöltek. Ezek az elemzőfában egy-egy csomópontot jelentenek.

### 5.1 Jellemzőkészlet

A tanító és a kiértékelő halmazon a jelöltekhez jellemzőket vettünk fel. Módszerünket gazdag jellemzőtérrel valósítottuk meg. Az eseménydetektálással kapcsolatos feladatokban gyakran használt jellemzőket mi is alkalmaztuk. Ezekon kívül újakkal is kibővítettük a jellemzőkészletünket. Az új jellemzőket a magyar szövegek tulajdonságai

<sup>1</sup> <http://www.cs.waikato.ac.nz/ml/weka/>

alapján választottuk ki. A jellemzőkhöz felhasználtuk a függőségi elemzőfát és a magyar WordNet-et is. A fő jellemző csoportokat több részre bontottunk, mivel a részek hatását külön-külön vizsgáljuk majd. Ezen csoportok közül a következők voltak az új, máshol ezen a területen nem látott jellemzőcsoportok: a két Morfológiai elemző együtt, az elemzőfa 1-2, szószák 1-3, WordNet jellemzők-2-4.

**Felszíni jellemzők:** *Bigramok, trigramok:* A vizsgált szavak végén lévő 2-es, 3-as betűcsoportok. *PositionInSentence:* a jelölt hányadik szó a mondatban. *NagyBetuNemMondatElejen:* Azok a nagybetűs szavak, amelyek nem a mondat elején állnak legtöbbször névelemek. Így ez a jellemző utalhat a nem-esemény jellegre.

**Morfológiai jellemzők-1:** Mivel a magyar nyelv igen gazdag morfológiával rendelkezik, ezért számos morfológia-alapú jellemzőt definiáltunk. Ebben a csoportban a Magyarlanc morfológiai elemzőjét használtuk fel. Jellemzőként definiáltuk az eseményjelöltek MSD morfológiai kódját, felhasználva a következő morfológiai jegyeket: *típus*(SubPos), *mód*(Mood), *eset*(Cas), *idő*(Tense), *személy*(PerP), *szám*(Num), *hatá-rozottság*(Def). További jellemzők: *Lemma:* a jelölt lemmája. *hasVerbRoot:* igéből képzett-e a jelölt. *SzofajElotte* és *SzofajUtana:* a jelölt előtti és utáni szó szófaja. *LegkozelebbiIgeMondatbanLemma:* a jelölthöz a mondatban legközelebb álló ige lemmája. *Igeto:* igéből képzett főnév esetén az alapige.

**Morfológiai jellemzők-2:** Ebben a csoportban a HunMorph morfológiai elemzőjét használtuk fel. *IgetoVan:* van-e igető. *IgebolFonevKepzo:* Igéből képzett főneveknél a képző. *IgeToHunMorph:* igéből képzett főnév esetén az alapige.

**Morfológiai jellemzők-3:** A Magyarlanc és a HunMorph morfológiai elemző is a többjelentésű szavak esetén minden jelentéshez megadnak külön morfológiai elemzést. Ebben a csoportban mindkét elemző esetén, egyértelműsítés nélkül, megadtuk a jelöltekhez minden jelentéshez tartozó ragokat, képzőket, jeleket.

**Elemzőfa jellemzők-1:** Ezeket a jellemzőket az függőségi elemzőfa alapján készítettük. *JeloltEdgeType:* A jelölt és az elemzőfában a felette levő szó közötti kapcsolat típusa. *JeloltEdgeTypeNE:* A jelölt NE (névelem) típussal kapcsolódik-e a felette levő szóhoz. Ez utal a jelölt nem esemény jellegére. *JeloltFelettLemmaFaban:* A jelölt feletti szó lemmája az elemzőfában. *JeloltFelettIgeLemmaFaban:* Az elemzőfában közvetlenül a jelölt felett lévő ige (ha van) lemmája. *KozvetlenSzintaktikaiKapcsolat:* Ha a jelölt fölött van közvetlenül ige az elemzőfában, akkor a kettő közötti szintaktikai kapcsolat típusa. *LegkozelebbiIgeFeletteFabanLemma, LegkozelebbiIgeFeletteTavolsagFaban:* Az elemzőfában jelölt feletti legközelebbi ige lemmája és annak távolsága a fában a jelölttől. *JeloltReszfaTokenekSzama:* Az elemzőfában a jelölt alá tartozó részfa elemeinek száma. *FeletteSzoEdgeType:* Az elemzőfában a jelölt feletti szó és az a feletti szó közötti kapcsolat típusa. A Magyarlanc elemzőnél az időhatározók, időbeliséget kifejező szavak az események felett helyezkednek el, ezért ezek jelenléte utalhat a jelölt eseményjellegére.

**Elemzőfa jellemzők-2:** Ha a jelölt nem közvetlenül kapcsolódik a felette levő igehez az elemzőfában, akkor részletesen jellemeztük a jelölt és az ige közötti útvonalat. *SzofajÚtvonal:* Egymás után írtuk a jelölt és az ige közötti csomópontok szófaját. Például: C↑S↑V↑C↑V↑V. *Lemmaútvonal:* Hasonlóan az előzőhöz a jelölt és az ige közötti lemmákat írtuk egymás után. Például: napoztatás↑és↑törölgetés↑hajszárító↑megszárít. *SzintaktikaiKapcsolat-Útvonal:* A jelölt és az

ige közötti útvonalon a szintaktikai kapcsolatok típusai egymás után. Például: OBL↑COORD↑SUBJ↑COORD↑CONJ↑.

**Szósák jellemzők-1:** Szósák modellt használtuk fel szócsoportok jellemzésére. *ReszfaLemmakSzoszakAtlag:* A jelölt alatt lévő részfa szavainak lemmáit reprezentáltuk szósák modellel. A tanító halmazon minden lemmához kiszámítottuk, hogy milyen arányban tartozott pozitív jelölt a részfájához. Majd minden jelölthöz kiszámítottuk a részfáját alkotó lemmákhoz tartozó arányok átlagát. Nagy átlag arra utal, hogy a jelölt részfájában fontos szavak vannak az eseményjelleg szempontjából. *ReszfaLemmakSzoszakLegnagyobb:* Hasonló az előzőhöz, de itt minden jelöltnél a részfájához tartozó lemmák közül azt választottuk ki, amelyikhez legnagyobb valószínűség tartozott. Nagy maximális valószínűség utal arra, hogy a jelölt részfájában van legalább egy olyan lemma, ami erősen fontos az eseményjelleg szempontjából. Ez a jellemző segít a részfa egy-egy fontos szavának felismerésében. *KozvetlenAlattaLemmakSzoszakAtlag* és *KozvetlenAlattaLemmakSzoszakLegnagyobb:* Az előzőkhöz hasonló, de itt nem a jelölt részfájához tartozó minden szót vizsgáltuk, hanem csak a részfa azon szavait, amelyek szintaktikailag kapcsolódnak a jelölthöz az elemzőfában. *KozvetlenAlattaEdgeTypeSzoszakAtlag* és *KozvetlenAlattaEdgeTypeSzoszakLegnagyobb:* Az előzőhöz hasonlóan, itt a jelölt és a hozzá szintaktikailag kapcsolódó szavak közötti kapcsolat típusát vizsgáltuk. *LemmaParseTreePathIgeigLemmakSzoszakAtlag* és *LemmaParseTreePathIgeigLemmakSzoszakLegnagyobb:* Ezeknél a szósákba a jelölt és az elemzőfában felette lévő ige közötti útvonalon található lemmák kerültek.

**Szósák jellemzők-2:** Ezekhez a jellemzőkhöz a lemmák a mondatból és nem az elemzőfából lettek kigyűjtve a szósákba. *MondatbanKornyezet-N-LemmakSzoszakAtlag* és *MondatbanKornyezet-N-LemmakSzoszakLegnagyobb:* A mondatban a jelölt N távolságú környezetét jellemeztük szósák modellel, N=3 és N=5 esetekben.

**WordNet jellemzők:** Ezekhez a jellemzőkhöz felhasználtuk a magyar WordNet [10] hiperním hierarchiájában található szemantikai kapcsolatokat. Mivel egy szóalakhoz több jelentés is tartozhat a WordNet-ben, ezért az egyes jelentések között egyértelműsítést végeztünk a Lesk algoritmussal[8]. A WordNetben a synsetekhez definíció és példamondatok tartoznak. Az algoritmus alapján, többjelentésű eseményjelölt esetén megszámoltuk, hogy az eseményjelölt szintaktikai környezetében lévő szavak közül hány található meg az egyes WordNet jelentések definíciói és példamondatai között. Azt a jelentést választottuk, amelyikkel a legtöbb volt közös szó.

**WordNet jellemzők-1:** *EsemenyszerusegekAlatt:* A magyar WordNetben van egy mesterséges synset, ami alá jellemzően események tartoznak. De vannak események ezen kívül is, és vannak ebben a csoportban olyanok is, amelyek nem igazi események. Ebben a jellemzőben megadtuk, hogy a jelölt beletartozik-e ezen synset hiponím hierarchiájába.

**WordNet jellemzők-2:** *WordNetSzoszakAtlag* és *WordNetSzoszakLegnagyobb:* A szósák jellemzőkhöz hasonlóan itt a szósákba a WordNetben a jelölt hiperním hierarchiájába tartozó szavakat vettük fel. *WordNetSzoszakLegnagyobbSynset:* Megadtuk a jelölt hiperním hierarchiájában lévő synset-ek közül azt, amelyik a legnagyobb arányban tartozik események hiperním hierarchiájába.

**WordNet jellemzők-3:** *WordNetHipernimSynsetekTanulobol* (bináris): Készítettünk egy halmazt, amibe kigyűjtöttük a tanító halmazból az esemény jelöltek hiperním hierarchiájának synset-jeit. Majd minden jelölthöz megadtuk, hogy a hiperním hierarchiájának synset-jei közül tartozik-e valamelyik ebbe a halmazba.

**WordNet jellemzők-4:** *WordNetLegjobbLemmakAlatt*: Kigyűjtöttük azokat a lemmákat, amelyek a tanító halmazon legnagyobb arányban voltak események. Majd a jelölteknel jelöltük, hogy a jelölt lemmája alatta van-e valamelyik ilyen kiemelt lemma hiponím hierarchiájának a WordNet-ben.

**Szósák jellemzők-3:** Először a *Szósák jellemzők 1-2* csoportoknál bemutatott minden esethez itt kiválasztottuk a legjobb elemeket a szósákokból 1-1 halmazba. Azokat, amelyek legnagyobb arányban tartoztak eseményekhez. Majd a következő jellemzőkkel jelöltük, hogy az adott jelölthöz tartozó szósák tartalmaz-e az adott halmaz elemei közül. *LegjobbWordNetSynsetek*: A jelölt hiperním hierarchiájába tartozó synsetek között van-e ami szerepel a LegjobbWordNetSynsetek halmazban. *LegjobbRészfaLemmak*: A jelölt részfaának lemmái között van-e olyan lemma, ami szerepel a LegjobbRészfaLemmak halmazban. *LegjobbLemmakÚtvonalIgeig*: A jelölt és az elemzőfában a legközelebbi ige közötti lemmák között van-e olyan lemma, ami szerepel a LegjobbÚtvonalLemmak halmazban. *LegjobbMondatbanKörnyezet-N-Lemmak*: A mondatban a jelölt N távolságú környezetében van-e olyan lemma, ami szerepel a LegjobbMondatbanKörnyezet-N-Lemmak között. Ezt megnéztük N=3 és N=5 esetekre is.

**Lista-jellemzők:** *FeletteLemmaldohatarozoListabol*: Listába kigyűjtöttünk gyakori idővel kapcsolatos szavakat. (például: előtt, folyamán) Ezek a szavak alatt az elemzőfában gyakran események vannak. Jellemzőként jelöltük, hogy a jelölt felett van-e ilyen idővel kapcsolatos kifejezés. *FeletteIgeAspektualisListabol*: Listába kigyűjtöttünk gyakori aspektuális igéket (például elkezd, folytatódik). Ezen igék alá tartozó főnevek gyakran események. Jelöltük, hogy a jelölt felett az elemzőfában van-e ilyen ige.

**Kombinált jellemzők-2 eleműek:** Ezeknél a jellemzőknél az előző jellemzők közül kombináltunk össze kettőt. *JeloltFelettLemmaFaban+JeloltEdge-Type*: Egy szó eseményjellegét gyakran pontosabban jelzi, ha a felette levő lemmát és a kettőjük közötti kapcsolatot együtt vizsgáljuk, mintha csak külön-külön vizsgálnánk azokat. Hasonlóan együtt vizsgáltuk a következőket:  
*JeloltFelettIgeLemmaFaban+JeloltEdgeTypeOBJ*,  
*JeloltFelettIgeLemmaFaban+JeloltEdgeTypeSUBJ*,  
*JeloltFelettLemmaFaban+LegjobbWordNetSynsetek*,  
*JeloltFelettIgeLemmaFaban+ LegjobbWordNetSynsetek*

**Kombinált jellemzők - 3 eleműek:** Az előző kételemű jellemzőkhöz hasonlóan itt három jellemzőt kombináltunk össze.  
*JeloltFelettLemmaFaban+EdgeType+WordNetLegjobbSynset*,  
*JeloltFelettIgeLemma-Faban+JeloltEdgeType+WordNetLegjobbSynset*,

## 5.2 További alkalmazott módszerek

A következő módszerek újaknak tekinthetők, mert ezeken a területeken nem láttuk máshol az alkalmazásukat. Mindegyik hasznos volt az eredménye alapján, így más NLP feladatoknál is hasznosak lehetnek.

*Statisztikai arány felhasználása az osztályozásnál.* A jelöltekhez a jellemzőket két módszer alapján választottuk ki. *Első módszernél* az előző részben bemutatott alapjellelmzőket használtuk fel. *Második módszernél* az alapjellelmzők helyett a tanító adaton számított statisztikai arányokat használtuk fel. A tanító halmaz alapján megszámloltuk minden jellemző esethez, hogy hány alkalommal fordult elő és ebből hányszor volt a jelölt *pozitív*. Ezek alapján kiszámítottuk a hozzá tartozó pozitív-arányt. Például ha a *Lemma* jellemzőnél a *Lemma = írás* eset 5-ször fordult elő és ebből 3-szor volt *pozitív eset*, akkor hozzá a 0,6-es pozitív-arány tartozott. Ebben az esetben az osztályozónak a jelöltekhez nem az alapjellelmzőt, hanem a hozzá tartozó arányt adtuk meg. Az előző példánál *Lemma-arány* = 0,6. Ezzel *jelentősen csökkentettük az osztályozó vektorterének méretét* az első módszerhez képest és így a *futási időt is*. Ez a kidolgozási időszakban hasznos volt. A két esetet összehasonlítva azt tapasztaltuk, hogy legtöbbször a *valószínűségi módszer* adta a legjobb eredményeket. És a futási idő is *70-80-szor gyorsabb*, mint az alapjellelmzők használata esetén.

*A jelöltek csoportokra bontása.* Az osztályozó hasonló tulajdonságú adathalmazon könnyebben találja meg a szabályokat, mint olyan halmazon, ami sokféle adatot tartalmaz. Ezért érdemes a jelölteket kisebb, hasonló tulajdonságú csoportokra bontani, (ezzel megkönnyíteni az osztályozó döntését). Majd a csoportok eredményeit a TP, TN, FP, FN eredmények alapján összegezni. Ennek megfelelően a jelöltjeinket két fő szempont szerint csoportosítottuk. Első lépésként a jelölteket két csoportra bontottuk: igéből képzett (deverbális) és nem igéből képzett (nem deverbális) főnevek. Hiszen e két csoport tagjai eltérően viselkednek. Az igéből képzett főnevek között sokkal nagyobb arányban vannak események. Másik csoportosítás a jelöltek lemmái alapján történt. Itt 3 alcsoportot képeztünk. Első csoportba azok a lemmák kerültek, amelyek nagy arányban események voltak a tanító halmazon. A másik csoportba a többi jelölt lemmája a tanító halmazról. Harmadik csoportba a kiértékelő halmazon azon jelöltek lemmái, amelyek nem szerepeltek a tanító halmazon. Így összesen  $2 \cdot 3 = 6$  csoportot képeztünk, és mindegyikre külön-külön végeztük el az osztályozást.

*Valószínűségi arányok felhasználása az osztályozás eredményeinek javítására.* Azokban az esetekben, amikor gyenge eredményt kaptunk, (általában a kis fedés miatt), akkor az osztályozás elvégzése után azon jelölteknél, amelyek a tanító halmazon nagy arányban voltak események, az értékelést a kiértékelő halmazon pozitívrá állítottuk.

Majd az eredményeknél látni fogjuk, hogy ezek a kiegészítő módszerek jelentősen javították az eredményeinket és a futási időt.



### 5.3 Jellemző-esetek számának csökkentése

A vektortér méretét csökkentettük a következő módszerrel: csak azokat a jellemző-előfordulásokat vettük fel az osztályozáshoz, amelyek a tanító halmazon legalább háromszor szerepeltek. Ezzel jelentősen csökkentettük a futási időt és csak az osztályozás szempontjából jelentéktelen jellemző-előfordulásokat hagytuk ki.

## 6 Eredmények

A kiértékelés során a pontosság (P), fedés (R) és F-mérték (F) metrikákat használtuk.

### 6.1 Baseline mérés

Modellünk hatékonyságának vizsgálatához Baseline mérést végeztünk. Ennek keretében a jelöltek közül az igei alapúakat vettük pozitív esetnek a többit pedig negatívnak. Ennek eredménye a következő volt: pontosság: 66,67, fedés: 47,57, F-mérték: 55,52. A további eredményeinken látni fogjuk, hogy gépi tanulási módszerünk jóval felülteljesítette a Baseline mérés eredményét.

### 6.2 Modellünk eredménye

Gépi tanulási módszerünk a következő eredményt érte el a teljes korpuszon az adott jellemzőkészlettel és a kiegészítő módszerekkel: Pontosság: 79,25, Fedés: 67,04, F-mérték: 71,94. A kiegészítő módszerek alkalmazása nélkül a következő eredményeket kaptuk: Pontosság: 70,32, Fedés: 60,51, F-mérték: 65,03. Látható, hogy a kiegészítő módszerekkel jelentős javulást tudunk elérni. A javulás 80%-át a jelöltek csoportosítása adta, a kisebb részt az osztályozás utáni javításból származott. Ha csak az első szempont szerint csoportosítottunk, akkor azt kaptuk, hogy az igéből képzett főnevek esetén a modell sokkal jobb eredményt ért el (F-mérték: 84,62), mint a nem igéből képzett főneveknél (F-mérték: 39,52).

Modellünket megvizsgáltuk az öt részkorpuszon is. Ezekre az 1. táblázatban látható F-mértékeket kaptuk.

**Table 1.** Eredmények a részkorpuszokon (%)

Részkorpusz	F-mérték
Szépirodalom-fogalmazás	75,24
Újsághírek	76,31
Üzleti rövidhírek	75,12
Számítógépes szövegek	71,57
Jogi szövegek	68,74

Legjobb eredményünket az újsághírek doménen, a legrosszabbat pedig a jogi szövegeken kaptuk.

### 6.3 Eredmények porlasztásos mérésrel

Megvizsgáltuk, hogy az egyes **jellemzőcsoportok** hogyan befolyásolják a gépi tanuló-rendszer eredményeit. Ehhez *porlasztásos mérést* végeztünk. Ekkor a teljes jellemzőkészletből elhagytuk az egyes jellemzőcsoportokat, majd a maradék jellemzőkre támaszkodva tanítottunk. Ennek eredményei a 2. táblázatban találhatóak. Az adatok azt mutatják, hogy az adott jellemzőcsoportot elhagyva hogyan változott az eredmény. A csökkenő (negatív) eredmény azt jelzi, hogy a vizsgált jellemzőcsoportnak pozitív hatása van az esemény felismerésben.

**Table 2.** A porlasztásos mérés eredményei (%)

Elhagyott jellemzők	Változás az F-mértékben
Felszíni jellemzők	-0,28
Morfológiai jellemzők-1	-2,51
Morfológiai jellemzők-2	-0,52
Morfológiai jellemzők-3	-2,01
Elemzőfa jellemzők-1	-1,92
Elemzőfa jellemzők-2	-0,52
Szósák jellemzők-1	-1,34
Szósák jellemzők-2	-2,42
Szósák jellemzők-3	-0,57
WordNet jellemzők-1	-0,32
WordNet jellemzők-2	-6,51
WordNet jellemzők-3	-0,53
WordNet jellemzők-4	-0,2
Lista jellemzők	0,0
Kombinált jellemzők - 2 eleműek	-0,79
Kombinált jellemzők - 3 eleműek	+0,1

Ha a hasonló jellemzőcsoportokat összevonjuk, akkor a következő eredményeket kapjuk az összevont csoportokra (3. táblázat):

**Table 3.** A porlasztásos mérés eredményei - összevonással (%)

Elhagyott jellemzők	Változás az F-mértékben
Morfológiai jellemzők	-1,63
Szósák jellemzők	-4,0
Elemzőfa jellemző	-1,56
WordNet jellemző	-7,7
Kombinált jellemzők	-0,95

A 2. és a 3. táblázat eredményein látszik, hogy majdnem minden jellemzőcsoportnak pozitív hatása volt a modell teljesítményére. Legjobb hatása a *WordNet* és a *Szó-*

zsák jellemzőknek volt, de sokat javítottak a *Morfológiai és az Elemzőfa* jellemzők is. Mindkét morfológiai elemző hatása pozitív volt. A WordNet jellemzők-2 részcsoporthoz volt a legjobb hatása (6.51%). Ebben használtuk együtt a WordNet-et a szózsák modellel. A Lista jellemzőknek nem volt hatása. Negatív hatása volt a 3 elemű kombinált jellemzőknek, de a 2 elemű kombinált jellemzők hasznosak voltak.

A modellünk, amelynek eredményét a 6.2-es fejezetben ismertettünk, már csak a pozitív hatású jellemzőket tartalmazta.

#### 6.4 Az eredmények összehasonlítása a kapcsolódó munkákkal.

Angol szövegekre Jeong és társa [6] 71,8%-os, Romeo és társai [12] 67%-os F-mértéket értek el. Olasz nyelvre Caselli [3] 69%-os, spanyol nyelvre Peris és társai [11] 59,6%-os F mértéket értek el. A kapcsolódó munkákkal összehasonlítva, eredményeink (F-mérték = 71,9%) jónak számítanak.

### Összegzés

Munkánkban bemutatunk gazdag jellemzőtérre alapuló gépi tanuló megközelítésünket, amely automatikusan képes magyar nyelvű szövegekben főnévi eseményeket detektálni. Öt részterületet vizsgáltunk meg, összesen 10 000 mondattal. Gazdag jellemzőtérre alapuló *jellemzőkészletünkben* felszíni, morfológiai, függőségi elemzőfa, szózsák, Wordnet, lista és kombinált jellemzőket használtunk fel. Ezen jellemzőcsoportok mellett kiegészítő módszereket is alkalmaztunk, amelyek javították modellünk hatékonyságát, valamint a futási időt. Algoritmusainkat tesztadatbázisokon kiértékelve, versenyképes eredményeket érnek el az eddig bemutatott angol és más nyelvű eredményekkel összehasonlítva.

### Bibliográfia

1. Bethard, S., Martin, J.H.: Identification of event mentions and their semantic class. In: Proceedings of the 2006 Conference on Empirical Methods in Natural Language Processing, pp. 146–154. Association for Computational Linguistics (2006)
2. Boguraev, B., Ando, R.K.: Effective use of Timebank for TimeML analysis. In: Schilder, F., Katz, G., Pustejovsky, J. (eds.) Annotating, Extracting and Reasoning about Time and Events. LNCS, vol. 4795, pp. 41–58. Springer, Heidelberg (2007)
3. Caselli, T., Russo, I., Rubino, F.: Recognizing deverbal events in context. In: Proceedings of CICLing 2011, poster session. Springer (2011)
4. Csendes, D., Csirik, J., Gyimóthy, T.: The Szeged corpus: a POS tagged and syntactically annotated hungarian natural language corpus. In: Sojka, P., Kopeček, I., Pala, K. (eds.) TSD 2004. LNCS (LNAI), vol. 3206, pp. 41–47. Springer, Heidelberg (2004)
5. Gorzitze, S., Pado, S.: Corpus-based acquisition of German event- and object denoting nouns. In: Proceedings of KONVENS 2012 (Main Track: Poster Presentations), pp. 259–263 (2012)
6. Jeong, Y., Myaeng, S.: Using syntactic dependencies and Wordnet classes for noun event recognition. In: The 2<sup>nd</sup> Workshop on Detection, Representation, and Exploitation of Events

- in the Semantic Web in Conjunction with the 11th International Semantic Web Conference, pp. 41–50 (2012)
7. Jurafsky, D., Martin, J.H.: *Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition*. Prentice-Hall, Upper Saddle River (2000)
  8. Lesk, M. Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In *SIGDOC '86: Proceedings of the 5th annual international conference on Systems documentation*, pp. 24–26, New York, NY, USA. ACM. (1986)
  9. Llorens, H., Saquete, E., Navarro-Colorado, B.: TimeML Events recognition and classification: learning CRF models with semantic roles. In: *Proceedings of the 23<sup>rd</sup> International Conference on Computational Linguistics*, pp. 725–733. Association for Computational Linguistics (2010)
  10. Miháلتz, M., Hatvani, Cs., Kuti, J., Szarvas, Gy., Csirik, J., Prószéký, G., Váradi, T.: Methods and results of the Hungarian WordNet project. In: Tanács, A., Csendes, D., Vincze, V., Fellbaum, C., Vossen, P., (eds.) *Proceedings of the Fourth Global WordNet Conference (GWC 2008)*, pp. 311–320. University of Szeged, Szeged (2008)
  11. Peris, A., Taule, M., Boleda, G., Rodríguez, H.: ADN-classifier: automatically assigning denotation types to nominalizations. In: *Proceedings of the Seventh LREC Conference*, 19–21 May 2010, Valetta, Malta, pp. 1422–1428 (2010)
  12. Romeo, L., Lebani, G.E., Bel, N., Lenci, A.: Choosing which to use? A study of distributional models for nominal lexical semantic classification. In: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC 2014)*, pp. 4366–4373 (2014)
  13. Sauri, R., Knippen, R., Verhagen, M., Pustejovsky, J.: Evita: a robust event recognizer for QA systems. In: *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing*, pp. 700–707. Association for Computational Linguistics (2005)
  14. Subecz, Z.: Detection and classification of events in Hungarian natural language texts. In: Sojka, P., Horak, A., Kopecek, I., Pala, K. (eds.) *TSD 2014. LNCS (LNAI)*, vol. 8655, pp. 68–75. Springer, Heidelberg (2014)
  15. Tron, V., Kornai, A., Gyepesi, G., Németh, L., Halácsy, P., Varga, D. Hunmorph: Open source word analysis. In *Proceedings of the Workshop on Software, Software '05*, pp. 77–85, Stroudsburg, PA, USA. Association for Computational Linguistics. (2005)
  16. Zsibrita, J., Vincze, V., Farkas, R.: magyarlanc: a toolkit for morphological and dependency parsing of Hungarian. In: *Proceedings of RANLP 2013*, pp. 763–771 (2013)